# PageX: An Integrated Document Processing and Management Software for Digital Library

**Hanchuan Peng, Zheru Chi, Wan-Chi Siu, and David D. Feng**

Department of Electronic & Information Eng.,
The Hong Kong Polytechnic University,
Hung Hom, Kowloon, Hong Kong.
Email: phc@eie.polyu.edu.hk

2000-Jan-10

# Abstract

*For digital libraries it is very important to design and to implement powerful engines to convert information on paper to electronic format. In this paper a PageX software is proposed as the integration of such engines, which include a set of intelligent document processing functions, a set of compact document management strategies, and a set of advanced accessories. With this software, a paper document will first be input as an optical image, which may be a mixture of graphics and text and may be skewed. The image will then be analyzed and decomposed into a series of component blocks, and encoded and stored in a structured and compact format. Well-developed accessory functions, including block editing and page annotation, page reconstruction and virtual editing, page matching and registration, document retrieval, etc., are provided to support advanced applications. With these carefully designed functions and strategies, PageX minimizes manual operations to a minimal degree.*

2000-Jan-10

# Main Process of Working


Image Acquiring Devices

**PageX**

**Binarization**
**Skew correction**
**Text blocks extraction**
**Script determination**
**Character segmentation**
**Thinning**
**Feature extraction**
……

**Page Format Database**

**PageX**

**Form Registration**

**PageX**

**Character Recognition**
**Page/Form Reformatting**
**Unrecognizable Block/Image Coding**
**Page Synthesis/Reconstruction**
……

# Main Structure of the System

## Optical Images



## Engine Set I: Page Analysis and Decomposition

(1) Foreground Extraction and Binarization
(2) Correction of Unknown Skew Angle
(3) Page Blocking
(4) Language Separation
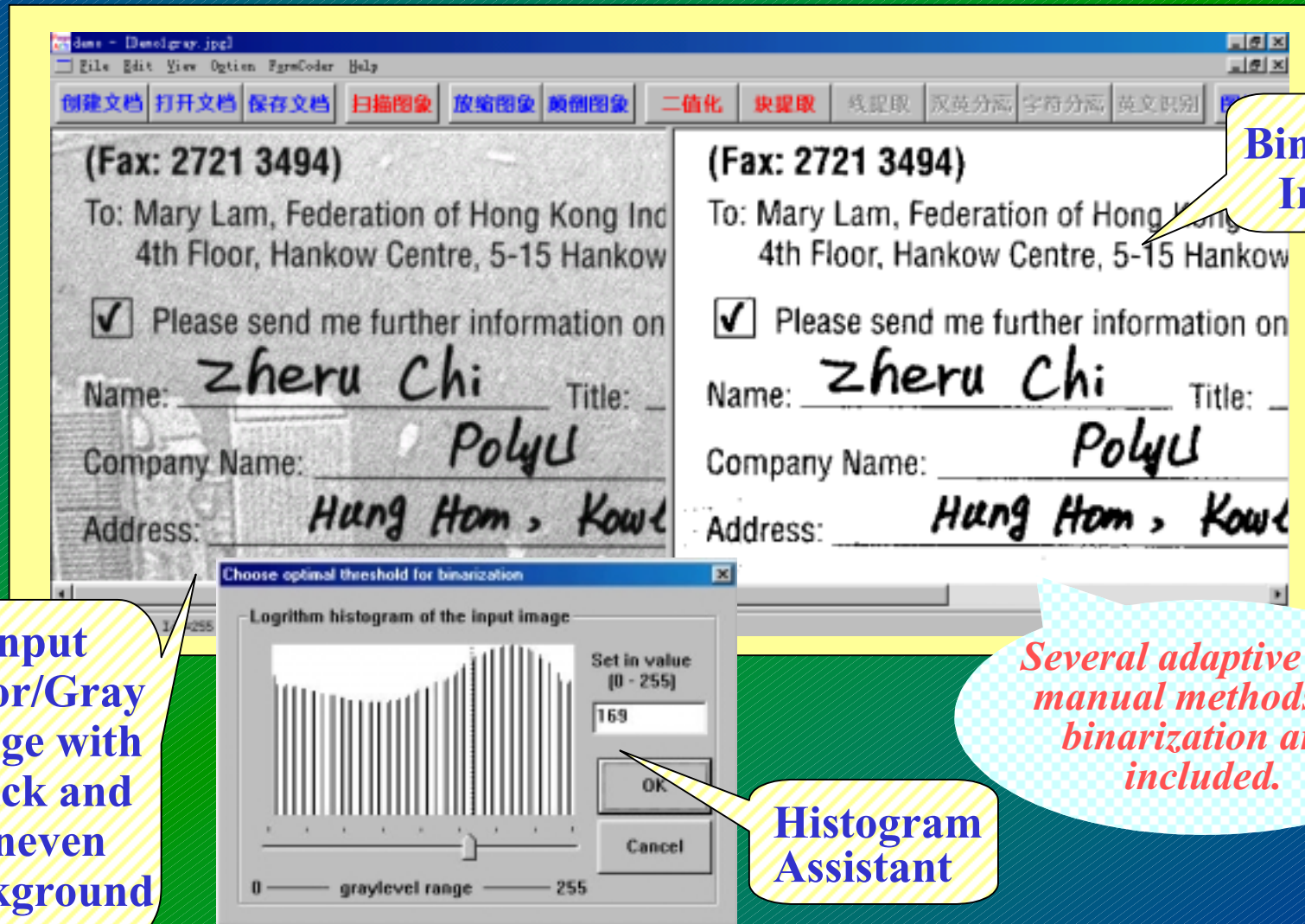(5) Character Separation
(6) Character Recognition
(7) ……

**Component Block List**

## Engine Set II: Block Coding and Page Management

(1) Block Compression and Decoding
(2) Page Layout Control
(3) Dynamically Adjustable Data Structures
(4) ……

**Compact e-Doc**

## Engine Set III: Advanced Accessories

(1) Virtual Editing
(2) Page Matching
(3) Database Linking
(4) Retrieval
(5) ……

**Compact e-Doc**

# Binarization



**Binarized Image**

**Input Color/Gray Image with Thick and Uneven Background**

**Histogram Assistant**

*Several adaptive and manual methods of binarization are included.*

# Skew Correction



**Input Image with Unknown Skew Angle**

**Output Image after automatic skew estimation and correction**

**Manual Skew Estimation Assistant**

*Both adaptive and manual skew correction methods are included.*

# Page Blocking

**Block Property Assistant**

**Block Registration Assistant**

**Block Property Assistant**

*Auto Page Blocking produces the basic components of the page image for further processing.*

# Script Determination

**Document image with multi-linguistic scripts**

Registered Design, Copyrig
can help to protect your new
shielding it from imitators.
外觀設計註冊　　版權及商
計不受抄襲

*Scripts of different languages can be distinguished and further encoded with different methods.*

Binarization ▶
Script determination ▶
Character recognition ▶
Thining ▶
Skew Correction ▶

✔ Momentum method (Peng)
Momentum method (Chi)
Neural net method (Zhu)
✔ Redo from the binary image
Just show current results

**Language Separation Assistant**

**English Parts**

Registered Design, Copyrig
can help to protect your new
shielding it from imitators.

**Chinese Parts**

外觀設計註冊　　版權及商標可
計不受抄襲

2000-Jan-10

# Virtual Editing



2000-Jan-10

# Page Matching

Component block list

· · · · · · Sequencing with size

Sequenced block list

· · · · · · $B$ $A$ · · · · · ·

Matching with position

Matching with size

Template block list

$B^T$

*Pages are parsed as blocks and the block sequences are matched.*